

# Kapitola 11.: Porovnání empirického a teoretického rozložení

## Cíl kapitoly

Po prostudování této kapitoly budete umět

- testovat hypotézu, že daný náhodný výběr pochází z rozložení s danou diskrétní či spojitou distribuční funkcí
- ověřovat podmínky dobré aproximace pro testy dobré shody
- pomocí jednoduchých testů testovat hypotézu, že daný náhodný výběr pochází z exponenciálního či Poissonova rozložení

## Časová zátěž

Na prostudování této kapitoly a splnění úkolů s ní spojených budete potřebovat asi 4 hodiny studia.

### 11.1. Motivace

Možnost použití statistických testů je podmíněna nějakými předpoklady o datech. Velmi často je to předpoklad o typu rozložení, z něhož získaná data pocházejí. Mnoho testů je založeno na předpokladu normality. (Testování normality bylo probráno ve 2. kapitole.) Opomíjení předpokladů o typu rozložení může v praxi vést i ke zcela zavádějícím výsledkům, proto je nutné věnovat tomuto problému patřičnou pozornost.

V této kapitole se seznámíme s testem dobré shody, který je (po splnění určitých předpokladů) použitelný k ověření shody empirického rozložení s jakýmkoliv teoretickým rozložením. Tato univerzálnost je ovšem provázena poněkud sníženou silou testu. Proto byly pro některá rozložení vyvinuty speciální testy využívající charakteristických vlastností těchto rozložení. Zde uvedeme tzv. jednoduché testy exponenciálního a Poissonova rozložení.

### 11.2. Testy dobré shody

#### 11.2.1. Popis testu

Testujeme hypotézu, která tvrdí, že náhodný výběr  $X_1, \dots, X_n$  pochází z rozložení s distribuční funkcí  $\Phi(x)$ .

- a) Je-li distribuční funkce spojitá, pak data rozdělíme do  $r$  třídících intervalů  $(u_j, u_{j+1})$ ,  $j = 1, \dots, r$ . Zjistíme absolutní četnost  $n_j$   $j$ -tého třídícího intervalu a vypočteme pravděpodobnost  $p_j$ , že náhodná veličina  $X$  s distribuční funkcí  $\Phi(x)$  se bude realizovat v  $j$ -tém třídícím intervalu. Platí-li nulová hypotéza, pak  $p_j = \Phi(u_{j+1}) - \Phi(u_j)$ .
- b) Má-li distribuční funkce nejvýše spočetně mnoho bodů nespojitosti, pak místo třídících intervalů použijeme varianty  $x_{[j]}$ ,  $j = 1, \dots, r$ . Pro variantu  $x_{[j]}$  zjistíme absolutní četnost  $n_j$  a vypočteme pravděpodobnost  $p_j$ , že náhodná veličina  $X$  s distribuční funkcí  $\Phi(x)$  se bude realizovat variantou  $x_{[j]}$ . Platí-li nulová hypotéza, pak  $p_j = \Phi(x_{[j]}) - \lim_{x \rightarrow x_{[j]}^-} \Phi(x) = P(X = x_{[j]})$ .

Testová statistika: 
$$K = \sum_{j=1}^r \frac{(n_j - np_j)^2}{np_j}$$
. Platí-li nulová hypotéza, pak  $K \approx \chi^2(r-1-p)$ , kde  $p$  je

počet odhadovaných parametrů daného rozložení. (Např. pro normální rozložení  $p = 2$ , protože z dat odhadujeme střední hodnotu a rozptyl.) Pokud žádný parametr nemusíme odhadovat, hovoříme o úplně specifikovaném problému. Nulovou hypotézu zamítáme na

asymptotické hladině významnosti  $\alpha$ , když  $K \geq \chi^2_{1-\alpha}(r-1-p)$ . Aproximace se považuje za vyhovující, když  $np_j \geq 5$ ,  $j = 1, \dots, r$ .

**Upozornění:** Při nesplnění podmínky  $np_j \geq 5$ ,  $j = 1, \dots, r$  je třeba některé intervaly resp. varianty slučovat, což vede ke ztrátě informace. Ve spojitém případě je hodnota testové statistiky  $K$  silně závislá na volbě třídících intervalů

**11.2.2. Příklad:** (Testování shody empirického a teoretického rozložení při úplně specifikovaném problému)

Ze souboru rodin s pěti dětmi bylo náhodně vybráno 84 rodin a byl zjišťován počet chlapců:

Počet chlapců	0	1	2	3	4	5
Počet rodin	3	10	22	31	14	4

Na asymptotické hladině významnosti 0,05 testujte hypotézu, že rozložení počtu chlapců se řídí binomickým rozložením  $Bi(5; 0,5)$ .

**Řešení:**

Pravděpodobnost, že náhodná veličina s rozložením  $Bi(5; 0,5)$  bude nabývat hodnot  $p_0, \dots, p_5$

je  $p_j = \binom{5}{j} \frac{1}{32}$ ,  $j=0,1,\dots,5$ .

Výpočty potřebné pro stanovení testové statistiky  $K$  uspořádáme do tabulky.

j	$n_j$	$p_j$	$np_j$
0	3	0,03125	84.0,03125=2,625
1	10	0,15625	84.0,15625=13,125
2	22	0,3125	84.0,3125=26,25
3	31	0,3125	84.0,3125=26,25
4	14	0,15625	84.0,15625=13,125
5	4	0,03125	84.0,03125=2,625

Podmínky dobré aproximace nejsou splněny, sloučíme tedy první dvě varianty a poslední dvě varianty.

j	$n_j$	$p_j$	$np_j$	$\frac{(n_j - np_j)^2}{np_j}$
0 a 1	13	0,1875	84.0,1875=15,75	0,480159
2	22	0,3125	84.0,3125=26,25	0,688095
3	31	0,3125	84.0,3125=26,25	0,859524
4 a 5	18	0,1875	84.0,1875=15,75	0,321429

Vypočteme realizaci testové statistiky:  $K = 0,48059 + 0,688095 + 0,859524 + 0,321429 = 2,3492$ , počet tříd  $r = 4$ , počet odhadovaných parametrů  $p = 0$ ,  $r - p - 1 = 3$ , kritický obor  $W = \langle \chi^2_{1-\alpha}(r-p-1), \infty \rangle = \langle \chi^2_{0,95}(3), \infty \rangle = \langle 7,8147; \infty \rangle$  Protože  $K \notin W$ , nulovou hypotézu nezamítáme na asymptotické hladině významnosti 0,05.

**Výpočet pomocí systému STATISTICA:**

Vytvoříme datový soubor se dvěma proměnnými a čtyřmi případy. Proměnná nj obsahuje zjištěné četnosti (po sloučení variant), proměnná npj pak teoretické četnosti.

Statistiky – Neparametrická statistika – Pozorované vs. očekávané  $\chi^2$  – OK – Proměnné – Pozorované četnosti nj, očekávané četnosti npj – OK – Výpočet.

		Pozorované vs. očekávané četnosti (Tabulka Chi-Kvadr. = 2,349206 sv = 3 p = ,503161)			
Případ		pozorov. nj	očekáv. npj	P - O	(P-O) <sup>2</sup> / O
C:	1	13,0000	15,7500	-2,7500	0,48015
C:	2	22,0000	26,2500	-4,2500	0,68809
C:	3	31,0000	26,2500	4,7500	0,85952
C:	4	18,0000	15,7500	2,2500	0,32142
Sčt		84,0000	84,0000	0,0000	2,34920

V záhlaví výstupní tabulky je uvedena hodnota testového kritéria (2,349206), počet stupňů volnosti = 3 a p-hodnota (0,503161). Nulová hypotéza se tedy nezamítá na asymptotické hladině významnosti 0,05.

**11.2.3. Příklad:** (Testování shody empirického a teoretického rozložení při neúplně specifikovaném problému – diskrétní případ)

V tabulce jsou rozříděny fotbalové zápasy určité soutěže podle počtu vstřelených branek.

Počet branek	0	1	2	3	4 a víc
Počet zápasů	19	30	17	10	8

Na asymptotické hladině významnosti 0,05 testujte hypotézu, že jde o výběr z Poissonova rozložení.

**Výpočet pomocí systému STATISTICA:**

Vytvoříme datový soubor s dvěma proměnnými a 5 případy. Proměnná POCET obsahuje počet vstřelených branek, proměnná CETNOST pak počet zápasů, v nichž bylo dosaženo zjištěného počtu branek.

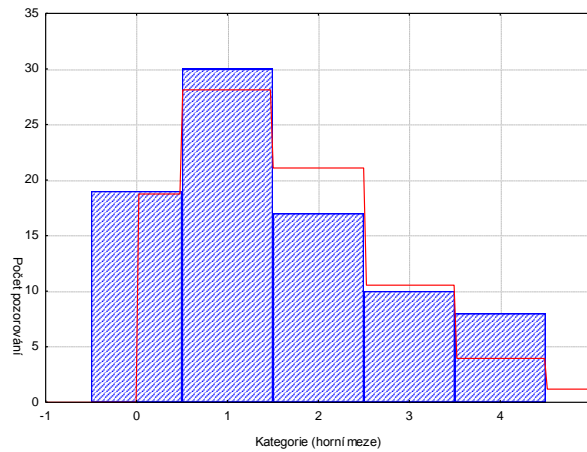
Statistiky – Prokládání rozdělení – Diskrétní rozdělení – Poissonovo – OK – Proměnná POCET – klikneme na ikonu se závažím – Proměnná vah CETNOST – Stav Zapnuto – OK – Výpočet.

Proměnná:POCET, Rozdělení:Poissonovo, Lambda = 1,500 (branky.sta) Chi-kvadrát = 2,07051, sv = 3, p = 0,55790								
Kategorie	Pozorované Četnosti	Kumulativ. Pozorované	Procent Pozorované	Kumul. % Pozorované	Očekáv. Četnosti	Kumulativ. Očekáv.	Procent Očekáv.	Kumul. % Očekáv.
<= 0,00000	19	19	22,6190	22,6190	18,7429	18,7429	22,3130	22,3130
1,00000	30	49	35,7142	58,3333	28,1144	46,8573	33,4695	55,7821
2,00000	17	66	20,2381	78,5714	21,0858	67,9431	25,1021	80,8842
3,00000	10	76	11,9047	90,4762	10,5429	78,4860	12,5510	93,4352
< Nekonečno	8	84	9,5238	100,0000	5,5139	84,0000	6,5642	100,0000

V tomto případě je parametr  $\lambda$  Poissonova rozložení neznámý, je odhadnut pomocí výběrového průměru a odhad činí 1,5. Podmínky dobré aproximace jsou splněny, dokonce všechny teoretické četnosti jsou větší než 5. Dále je v záhlaví výstupní tabulky uvedena hodnota testového kritéria (2,07051), počet stupňů volnosti  $r - p - 1 = 5 - 1 - 1 = 3$  a p-

hodnota (0,5578). Nulová hypotéza se tedy nezamítá na asymptotické hladině významnosti 0,05.

Pro vytvoření grafu se vrátíme do Proložení diskretních rozložení – Základní výsledky – Graf pozorovaného a očekávaného rozdělení.



#### 11.2.4. Příklad: (Testování shody empirického a teoretického rozložení při neúplně specifikovaném problému – spojitý případ)

U 48 studentek VŠE v Praze byla zjišťována výška (v cm):

165	170	170	179	170	168	174	162	167	165	170	173	183	176	165	168
171	178	168	168	169	163	172	184	176	175	176	169	168	170	166	160
167	162	162	166	170	168	155	162	169	166	160	169	165	163	168	163

Pomocí testu dobré shody testujte na hladině významnosti 0,05 hypotézu, že data pocházejí z normálního rozložení. Pomocí N-P grafu posuďte vizuálně předpoklad normality.

#### Výpočet pomocí systému STATISTICA:

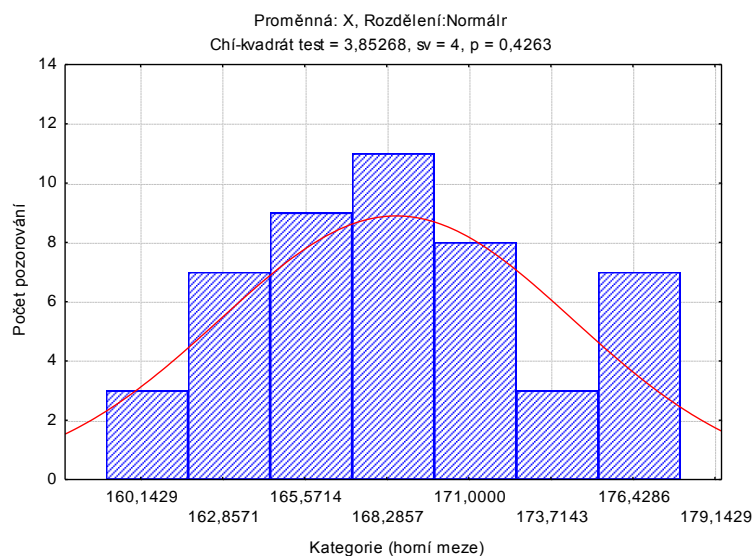
Statistiky - Prokládání rozdělení – ponecháme implicitní nastavení na normální rozložení – OK – Proměnná X – OK – na záložce Parametry změníme Počet kategorií na 7 (podle Sturgesova pravidla) – Výpočet.

Proměnná: X, Rozdělení: Normální (výška.sta) Chí-kvadrát = 1,09280, sv = 1 (uprav.), p = 0,29585								
Horní hranice	Pozorované četnosti	Kumulativ. Pozorované	Procent Pozorované	Kumul. % Pozorované	Očekáv. četnosti	Kumulativ. Očekáv.	Procent Očekáv.	Kumul. % Očekáv.
<= 157,14286	1	1	2,0833	2,0833	1,1970	1,1970	2,4938	2,4938
162,28571	6	7	12,5000	14,5833	5,5148	6,7118	11,4892	13,9833
167,42857	12	19	25,0000	39,5833	13,4622	20,1740	28,0462	42,0296
172,57143	19	38	39,5833	79,1667	15,8914	36,0655	33,1072	75,1366
177,71429	6	44	12,5000	91,6667	9,0770	45,1425	18,9104	94,0476
182,85714	2	46	4,1666	95,8333	2,5036	47,6462	5,2159	99,2622
< Nekonečno	2	48	4,1666	100,0000	0,3538	48,0000	0,7370	100,0000

Při tomto rozřídění dat do 7 intervalů nejsou splněny podmínky dobré aproximace, ve třech intervalech jsou teoretické četnosti pod 5. Změníme tedy dolní mez na 159 a horní na 178.

Proměnná: X, Rozdělení: Normální (vyska.sta) Chi-kvadrát = 3,85268, sv = 4, p = 0,42631								
Homí hranice	Pozorované Četnosti	Kumulativ. Pozorované	Procent Pozorované	Kumul. % Pozorované	Očekáv. Četnosti	Kumulativ. Očekáv.	Procent Očekáv.	Kumul. % Očekáv.
<= 161,71429	3	3	6,2500	6,2500	5,72299	5,7230	11,9229	11,9229
164,42857	7	10	14,5833	20,8333	5,67594	11,3989	11,8248	23,7477
167,14286	9	19	18,7500	39,5833	7,86263	19,2615	16,3804	40,1281
169,85714	11	30	22,9167	62,5000	8,81245	28,0740	18,3592	58,4873
172,57143	8	38	16,6667	79,1667	7,99151	36,0655	16,6489	75,1362
175,28571	3	41	6,2500	85,4167	5,86355	41,9291	12,2157	87,3520
< Nekonečno	7	48	14,5833	100,0000	6,07089	48,0000	12,6477	100,0000

V tomto případě jsou podmínky dobré aproximace splněny. Testová statistika se realizuje hodnotou 3,85268, p-hodnota je 0,42631, tedy na asymptotické hladině významnosti 0,05 hypotézu o normalitě nezamítáme. Podívejme se ještě na histogram s proloženou Gaussovou křivkou: Na záložce Základní výsledky zvolíme Graf pozorovaného a očekávaného rozdělení.



### 11.3. Jednoduchý test exponenciálního a Poissonova rozložení

#### 11.3.1. Jednoduchý test exponenciálního rozložení

Testujeme hypotézu, která tvrdí, že náhodný výběr  $X_1, \dots, X_n$  pochází z exponenciálního rozložení. Označme  $M$  výběrový průměr a  $S^2$  výběrový rozptyl tohoto náhodného výběru. Víme, že střední hodnota náhodné veličiny  $X \sim \text{Ex}(\lambda)$  je  $E(X) = 1/\lambda$  a rozptyl je  $D(X) = 1/\lambda^2$ .

Test založíme na statistice  $K = \frac{(n-1)S^2}{M^2}$ , která se v případě platnosti  $H_0$  asymptoticky řídí

rozložením  $\chi^2(n-1)$ . Kritický obor:  $W = \langle 0, \chi^2_{\alpha/2}(n-1) \rangle \cup \langle \chi^2_{1-\alpha/2}(n-1), \infty \rangle$ . Jestliže  $K \in W$ ,  $H_0$  zamítáme na asymptotické hladině významnosti  $\alpha$ .

#### 11.3.2. Příklad

Byla zkoumána doba životnosti 45 součástek (v hodinách). Průměrná životnost byla  $m = 99,93$  a rozptyl  $s^2 = 7328,91$ . Na asymptotické hladině významnosti 0,05 testujte hypotézu, že daný náhodný výběr pochází z exponenciálního rozložení.

**Řešení:**

Testovou statistiku  $K$  vypočteme podle vzorce  $K = \frac{(n-1)S^2}{M^2}$ . Kritický obor má tvar:

$$W = \langle 0; \chi^2_{\alpha/2}(n-1) \rangle \cup \langle \chi^2_{1-\alpha/2}(n-1); \infty \rangle. \text{ V našem případě } K = 32,2924,$$

$W = \langle 0; 27,575 \rangle \cup \langle 64,202; \infty \rangle$ ,  $H_0$  tedy nezamítáme na asymptotické hladině významnosti 0,05.

### 11.3.3. Jednoduchý test Poissonova rozložení

Testujeme hypotézu, která tvrdí, že náhodný výběr  $X_1, \dots, X_n$  pochází z Poissonova rozložení. Označme  $M$  výběrový průměr a  $S^2$  výběrový rozptyl tohoto náhodného výběru. Víme, že střední hodnota náhodné veličiny  $X \sim \text{Po}(\lambda)$  je  $E(X) = \lambda$  a rozptyl je  $D(X) = \lambda$ . Test založíme

na statistice  $K = \frac{(n-1)S^2}{M}$ , která se v případě platnosti  $H_0$  asymptoticky řídí rozložením

$\chi^2(n-1)$ . Kritický obor:  $W = \langle 0, \chi^2_{\alpha/2}(n-1) \rangle \cup \langle \chi^2_{1-\alpha/2}(n-1), \infty \rangle$ . Jestliže  $K \in W$ ,  $H_0$  zamítáme na asymptotické hladině významnosti  $\alpha$ .

### 11.3.4. Příklad

Studujeme rozložení počtu pacientů, kteří během 75 dnů přijdou na pohotovost.

Osmihodinovou pracovní dobu rozdělíme do půlhodinových intervalů a v každém intervalu zjistíme počet příchozích pacientů:

Počet pacientů	Pozorovaná četnost
0	79
1	188
2	282
3	275
4	196
5	114
6	45
7	10
8	7
9	3
10	1

Na asymptotické hladině významnosti 0,05 testujte hypotézu, že daný náhodný výběr pochází z Poissonova rozložení.

#### Řešení:

Celkový počet pacientů je  $n = 1200$ . Realizaci výběrového průměru  $M$  získáme jako vážený průměr počtu pacientů ( $m = 2,8033$ ) a realizaci výběrového rozptylu  $S^2$  získáme jako vážený rozptyl počtu pacientů ( $s^2 = 2,7086$ ). Testovou statistiku vypočteme podle

vzorce  $K = \frac{(n-1)S^2}{M}$ , tedy  $K = 1158,5$ , kritický obor

$$W = \langle 0, \chi^2_{\alpha/2}(n-1) \rangle \cup \langle \chi^2_{1-\alpha/2}(n-1), \infty \rangle = \langle 0, \chi^2_{0,025}(1199) \rangle \cup \langle \chi^2_{0,975}(1199), \infty \rangle = \\ = \langle 0; 1104,93 \rangle \cup \langle 1296,86; \infty \rangle$$

Protože testová statistika se nerealizuje v kritickém oboru,  $H_0$  nezamítáme na asymptotické hladině významnosti 0,05.

## Shrnutí

K ověření shody empirického rozložení s teoretickým rozložením se používají různé metody. Zvláštní postavení mezi nimi zaujímají metody zaměřené na ověřování normality dat. S nimi jsme se seznámili ve 2. kapitole.

Obecně je na ověření předpokladu o typu rozložení, z něhož pochází daný náhodný výběr, určen *chi-kvadrát test dobré shody*. Ten je založen na porovnání empirických četností jednotlivých variant či třídících intervalů s tzv. *teoretickými četnostmi*. Velké odchylky mezi empirickými a teoretickými četnostmi vedou k velkým hodnotám testového kritéria a tudíž k zamítnutí nulové hypotézy. Test dobré shody lze aplikovat pouze při splnění *předpokladů dobré aproximace*.

Pro ověřování shody empirických dat s exponenciálním či Poissonovým rozložením byly vyvinuty *jednoduché testy*, které využívají pouze znalosti rozsahu výběru, výběrového průměru a výběrového rozptylu.

## Kontrolní otázky

1. Popište provedení testu dobré shody pro náhodný výběr z diskrétního rozložení a pro náhodný výběr ze spojitého rozložení.
2. Jakým rozložením se asymptoticky řídí testová statistika testu dobré shody v případě platnosti nulové hypotézy?
3. Za jakých podmínek lze použít test dobré shody?
4. Popište jednoduchý test exponenciálního rozložení a Poissonova rozložení.

## Autokorekční test

1. Při 600 hodech kostkou byly zjištěny tyto četnosti: 85 x jednička, 99 x dvojka, 91 x trojka, 108 x čtyřka, 119 x pětka, 98 x šestka. Příspěvek šestky do testové statistiky  $K$  je

- a) 0,04
- b) 0
- c) 4

2. Uvažme zadání z otázky 1. Pokud je pravdivá hypotéza, že kostka je homogenní, pak testová statistika  $K$  se asymptoticky řídí *chi-kvadrát rozložením* s počtem stupňů volnosti

- a) 4
- b) 6
- c) 5

3. Na základě náhodného výběru rozsahu 537 z diskrétního rozložení je na asymptotické hladině významnosti 0,01 testem dobré shody ověřována hypotéza, že tento výběr pochází z Poissonova rozložení, přičemž parametr  $\lambda$  není znám. V datech se vyskytuje 5 variant náhodné veličiny  $X$ . Kritický obor pro test nulové hypotézy má tvar:

- a)  $W = \langle 0; 0,072 \rangle \cup \langle 12,838; \infty \rangle$
- b)  $W = \langle 12,838; \infty \rangle$
- c)  $W = \langle 9,348; \infty \rangle$

4. Jednoduchým testem provedeným na hladině významnosti 0,05 chceme ověřit hypotézu, že náhodný výběr rozsahu 43 pochází z exponenciálního rozložení, přičemž výběrový průměr

nabyl hodnoty 20,2558 a výběrová směrodatná odchylka 22,5051. Testová statistika se realizuje hodnotou

- a) 51,8457
- b) 2,3037
- c) 46,664

Správné odpovědi: 1a) 2c) 3b) 4a)

## Příklady

1. Ve svých pokusech pozoroval J. G. Mendel 10 rostlin hrachu a na každé z nich počet žlutých a zelených semen. Výsledky pokusu:

č. rostliny	1	2	3	4	5	6	7	8	9	10
počet žlutých	25	32	14	70	24	20	32	44	50	44
počet zelených	11	7	5	27	13	6	13	9	14	18
celkem	36	39	19	97	37	26	45	53	64	62

Z genetických modelů vyplývá, že pravděpodobnost výskytu žlutého semene by měla být 0,75 a zeleného 0,25. Na asymptotické hladině významnosti 0,05 testujte hypotézu, že výsledky Mendelových pokusů se shodují s modelem.

Výsledek:

Testová statistika  $K = 1,797495$ , kritický obor  $W = \langle \chi^2_{0,95}(9), \infty \rangle = \langle 16,9; \infty \rangle$ , nulovou hypotézu nezamítáme na asymptotické hladině významnosti 0,05.

2. Při 60 hodech kostkou jsme dosáhli těchto výsledků: 9 x jednička, 11 x dvojka, 10 x trojka, 13 x čtyřka, 11 x pětka a 6 x šestka. Na asymptotické hladině významnosti 0,05 testujte hypotézu, že kostka je homogenní.

Výsledek:

Testová statistika  $K = 2,8$ , kritický obor  $W = \langle \chi^2_{0,95}(5), \infty \rangle = \langle 11,07; \infty \rangle$ , nulovou hypotézu nezamítáme na asymptotické hladině významnosti 0,05.

3. Ze záznamů autosalónu byl ve 100 náhodně vybraných dnech zjištěn počet prodaných aut.

Počet prodaných aut za den	0	1	2	3	4	5 a víc
Počet dnů	9	43	29	11	5	3

Na asymptotické hladině významnosti 0,05 testujte hypotézu, že počet prodaných aut za den se řídí Poissonovým rozložením.

Výsledek:

Odhad parametru  $\lambda$  získaný pomocí výběrového průměru je 1,7.

Testová statistika  $K = 10,8891$ , kritický obor  $W = \langle \chi^2_{0,95}(4), \infty \rangle = \langle 9,488; \infty \rangle$ , nulovou hypotézu zamítáme na asymptotické hladině významnosti 0,05.

4. Při parlamentních volbách získaly 4 nejsilnější strany 30%, 20%, 15% a 10% hlasů, zbytek hlasů byl rozdělen mezi ostatní strany. Při volbách do obecního zastupitelstva v jedné obci získaly zmíněné strany (ve stejném pořadí) 1400, 900, 900 a 600 hlasů z 5000 odevzdaných



hlasů. Na asymptotické hladině významnosti 0,05 testujte hypotézu, že rozložení hlasů při parlamentních a místních volbách (v uvedené obci) je stejné.

Výsledek:

Testová statistika  $K = 68,67$ , kritický obor  $W = \langle \chi^2_{0,95}(4), \infty \rangle = \langle 9,488; \infty \rangle$ , nulovou hypotézu zamítáme na asymptotické hladině významnosti 0,05. S rizikem omylu nejvýše 5% jsme prokázali, že rozložení hlasů při parlamentních volbách a volbách v uvedené obci se liší.